

# Fully online clustering of evolving data streams into arbitrarily shaped clusters

Hyde, Richard; Angelov, Plamen; MacKenzie, A. R.

DOI:

[10.1016/j.ins.2016.12.004](https://doi.org/10.1016/j.ins.2016.12.004)

License:

Creative Commons: Attribution-NonCommercial-NoDerivs (CC BY-NC-ND)

*Document Version*

Peer reviewed version

*Citation for published version (Harvard):*

Hyde, R, Angelov, P & MacKenzie, AR 2017, 'Fully online clustering of evolving data streams into arbitrarily shaped clusters', *Information Sciences*, vol. 382-383, pp. 96-114. <https://doi.org/10.1016/j.ins.2016.12.004>

[Link to publication on Research at Birmingham portal](#)

## General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

## Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

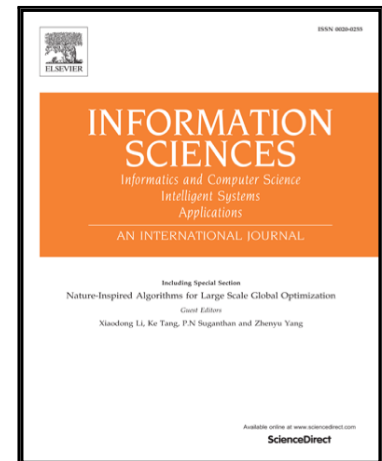
If you believe that this is the case for this document, please contact [UBIRA@lists.bham.ac.uk](mailto:UBIRA@lists.bham.ac.uk) providing details and we will remove access to the work immediately and investigate.

# Accepted Manuscript

## Fully Online Clustering of Evolving Data Streams into Arbitrarily Shaped Clusters

Richard Hyde, Plamen Angelov, A.R. MacKenzie

PII: S0020-0255(16)31924-7  
DOI: [10.1016/j.ins.2016.12.004](https://doi.org/10.1016/j.ins.2016.12.004)  
Reference: INS 12639



To appear in: *Information Sciences*

Received date: 3 February 2016  
Revised date: 25 November 2016  
Accepted date: 4 December 2016

Please cite this article as: Richard Hyde, Plamen Angelov, A.R. MacKenzie, Fully Online Clustering of Evolving Data Streams into Arbitrarily Shaped Clusters, *Information Sciences* (2016), doi: [10.1016/j.ins.2016.12.004](https://doi.org/10.1016/j.ins.2016.12.004)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# Fully Online Clustering of Evolving Data Streams into Arbitrarily Shaped Clusters

Richard Hyde, Plamen Angelov

*Data Science Group, School of Computing and Communications, Lancaster University,  
Lancaster, LA1 4WA, UK*

A R MacKenzie

*Birmingham Institute of Forest Research (BIFoR), University of Birmingham, Edgbaston,  
Birmingham, B15 2TT, UK*

---

## Abstract

In recent times there has been an increase in data availability in continuous data streams and clustering of this data has many advantages in data analysis. It is often the case that these data streams are not stationary, but evolve over time, and also that the clusters are not regular shapes but form arbitrary shapes in the data space. Previous techniques for clustering such data streams are either hybrid online / offline methods, windowed offline methods, or find only hyper-elliptical clusters. In this paper we present a fully online technique for clustering evolving data streams into arbitrary shaped clusters. It is a two stage technique that is accurate, robust to noise, computationally and memory efficient, with a low time penalty as the number of data dimensions increases. The first stage of the technique produces micro-clusters and the second stage combines these micro-clusters into macro-clusters. Dimensional stability and high speed is achieved through keeping the calculations both simple and minimal using hyper-spherical micro-clusters. By maintaining a graph structure, where the micro-clusters are the nodes and the edges are its pairs with intersecting micro-clusters, we minimise the calculations required for macro-cluster maintenance. The micro-clusters themselves are described in such a way that there is no calculation required for the core and shell regions and no separate definition of outer micro-clusters necessary. We demonstrate the ability of the proposed technique to join and separate macro-clusters as they evolve in a fully online manner. There are no other fully online techniques that the authors are aware of and so we compare the technique with popular online / offline hybrid alternatives for accuracy, purity and speed. The technique is then applied to real atmospheric science data streams and used to discover short term, long term and seasonal drift and their effects on anomaly detection. As well as having favourable computational characteristics, the technique can add analytic value over hyper-elliptical methods by characterising the cluster hyper-shape using Euclidean or fractal shape factors. Because the technique records macro-

clusters as graphs, further analytic value accrues from characterising the order, degree, and completeness of the cluster-graphs as they evolve over time.

*Keywords:* online, evolving, clustering, micro-cluster, arbitrary shape

## 1. Introduction

Recent technological advances in many disciplines have seen an increase in the amount of data being provided in continuous streams of data, i.e. ‘online data’. These data streams range from machine condition monitoring and atmospheric science data to social media analysis. The analysis and clustering of data streams has become increasingly important [3]. However, condition monitoring can suffer from sensor drift due to ageing, temperature fluctuations, modifications or upgrades to machine components, changes in load or type of use. Environmental monitoring will also be affected by sensor drift, but also seasonal variations and secular trends due to technological, socio-economic or climate change. While seasonality and other cyclic periodicities can be moved relatively easily off-line, any attempt to do this online renders the analysis vulnerable to aliasing changing seasonal cycles into secular changes. Other problem datasets are short-term but high-dimension and rapidly changing: chemical batch processors [11], environmental mesocosms [29], or ecological manipulation experiments [23], for instance. Social media analysis will be affected by the inevitable changes in peoples’ taste, population changes and many other influences. In examples such as these the assumption of a stationary data environment is invalid and techniques for data analysis need to be capable of coping with the evolution of these data streams. It is often the case in such data, particularly that incorporating spatial or relational information, that clusters of related data will not be hyper-elliptical and will fall into arbitrarily shaped groupings. The cases for arbitrary shaped clusters are well established and found in many sources [8, 26, 24]. Specifically a case such as that shown in [3] demonstrates the need for evolving clusters of arbitrary shapes - as the nature of the landscape changes over time, so must the clusters.

The ability to adapt our analytic to these secular (non-periodic) changes requires not only a method of reducing the importance of old data but also a way to divide previously singular clusters of data into multiple clusters. With the previously available techniques discussed in section 2 this is achieved, not by dividing the clusters in an online manner, but rather by re-clustering using an offline clustering technique on demand. With ever-increasing data sets, i.e. ‘Big Data’, the need to discard or archive the data after processing once becomes necessary for both computational and memory efficiency.

The technique presented in this paper has two distinct stages. The first adds data to current micro-clusters and adjust their information, or creates micro-clusters when data samples occur in un-clustered data space. The radius of the micro-clusters,  $r_0$ , is fixed and should be as small as is practical. In this newly proposed method we use a simple linear ageing process which reduces the ‘Energy’ of a micro-cluster and allows unused micro-clusters to die out completely.

Alternative ageing techniques could be used including those exponential types that leave micro-clusters present, with insignificant Energy, but allows them to be ‘re-born’ and become relevant in the future with further data. The micro-cluster Energy is renewed every time they receive new data. When no data is received the micro-clusters lose some Energy, gradually fading out. If no data is received for a long time the micro-cluster Energy will reach zero and they are discarded.

The second stage searches for overlapping micro-clusters. The micro-clusters are defined as having a kernel region  $\leq 0.5r_0$  and a shell region  $> 0.5r_0$ . By only connecting those micro-clusters whose kernel regions overlap into another micro-cluster shell we automatically determine edge micro-clusters. Micro-clusters which do not have at least the user-specified local density, i.e. the minimum number of samples within the radius, remain as separate outlier micro-clusters. Each macro-cluster consists of the graph of intersecting micro-clusters; the adjacency relations for each micro-cluster are stored as a property of that micro-cluster. For convenience, we call micro-clusters in adjacency relations (i.e. intersecting micro-clusters) Edges. Those micro-clusters with no Edges define graphs of order 1 without edges (i.e. intersections) and constitute a macro-cluster graph by themselves. Using this graph structure reduces the calculations required to separate clusters if a cluster dies and breaks a chain graph resulting in two groups of micro-clusters no longer being connected.

We call this technique Clustering of Evolving Data-streams into Arbitrary Shapes (CEDAS). We demonstrate the efficacy of CEDAS by testing the algorithm for speed and dimensional effects on synthetic data sets. Application of CEDAS to the KDDCup99 data stream set is used to compare cluster purity with DenStream and MR-Stream and also to demonstrate the ability of CEDAS to deal with Big Data, adapt to evolving data streams and detect internet intrusion attacks with high accuracy. We then apply the algorithm to real-world London Air Quality [11] atmospheric monitoring data to demonstrate how, by varying the micro-cluster decay time, we can differentiate between short term and long term secular change to discover temporally local anomalies and extremes of the overall distribution.

The rest of the paper is structured as follows: Section 2 provides a review of the current state of the art. Section 3 describes the principles and methodology behind the CEDAS algorithm and provides a description of the pseudo-code. Section 4 describes the data sets and the methodology of their use throughout the analysis parts of the paper and provides analysis of the performance of the proposed algorithm and comparisons to alternative techniques. The findings are summarised in section 5 and finally we consider some directions for future work in Section 6.

## 2. State of the Art

We refer to hyper-ellipsoidal techniques where the cluster membership is based on a distance based measure from a cluster centre. Where these clusters may overlap we acknowledge that they typically form centroidal veronoi

85 tessellations. Clustering techniques such as ELM [10], DEC [4] provide online  
clustering of data streams. Both of these techniques operate on data streams  
and provide clustering results online but are limited to hyper-ellipsoidal cluster  
shapes. The basis for ELM is to store the local mean as a cluster centre and to  
adjust the cluster centre and radii as more data arrives. DEC maintains a list of  
90 core and non-core clusters defined by the weight of the cluster. The weight de-  
cays over time or is increased as new data samples join the cluster. In this way,  
core clusters may decay to non-core, non-core clusters may disappear or increase  
their weight to become core clusters or new, non-core, clusters may be created.  
In both techniques the clusters that are created are hyper-ellipsoidal. In the  
95 case of concave cluster shapes DEC may create many smaller hyper-ellipsoidal  
clusters or one large cluster encapsulating all the data depending on the user  
parameter values.

Other existing clustering methods such as Chameleon [19], DBScan [12] and  
SPARCL [9] are all techniques for clustering arbitrary shapes offline. Sparcl  
100 utilises a two-layer approach whereby k-means [22] clustering is used to create  
a large number of micro-cluster centres in the first layer. These micro-cluster  
centres are then further clustered using a hierarchical approach. Chameleon  
and DBScan are techniques that successfully cluster arbitrary shapes, however,  
both work offline and therefore require the full data set. An incremental version  
105 of DBScan [7] was proposed which allows for incremental modification of the  
clusters. However, after each increment the micro-clusters are re-built and so  
require the full data from each increment to be available.

A method known as DenStream was proposed in [7] based on the CluStream  
[1]. A set of core- and potential-micro-clusters are maintained online. Each  
110 micro-cluster is created from a stored set of data with a decaying weight. By  
decaying the data samples those with a weight below a user-specified threshold  
are discarded and the memory requirement is limited somewhat although this  
loses potentially useful micro-clusters. The technique has an initialisation phase,  
using DBScan, to create an initial set of micro-clusters. Additionally, while the  
115 micro-clusters are maintained in an online fashion the process of combining the  
micro-clusters into final clusters is an off-line approach carried out on demand.  
DenStream is capable of finding arbitrarily shaped clusters as its 2nd stage  
clustering is based on DBScan whereas CluStream, with its 2nd stage based on  
k-means [15], finds hyper-elliptical clusters.

120 Two developments of DenStream known as SDStream [27] and rDenStream  
[20] improve on the basic DenStream algorithm. SDStream is based on SWClus-  
tering Algorithm described in [30] and is an offline approach, repeated tempo-  
rally at incremental time windows. The authors claim improved quality of  
clusters over CluStream [1] however, it remains an incremental offline approach  
requiring storage of past data. rDenstream is a three stage clustering technique  
125 also based on DenStream. In rDenStream any discarded clusters are retained  
in memory and may be re-introduced to improve the clustering at a later time.  
Processing of this discarded data requires additional processing time and mem-  
ory allocation. The 2nd stage of these techniques uses DBScan which has been  
130 demonstrated have an order of  $D^4$  time penalty, where  $D$  is the number of di-

mensions [18]. DBScan becomes impractical for big data and high dimensions and as a result techniques based around DBScan may also suffer from this limitation. The combined online-offline approaches limit this time penalty by only applying the offline DBScan macro-clustering function ‘on demand’ at reduced frequency.

Grid-based technique MR-Stream [28], divides the data space into a tree of grids of decreasing size (increasing resolution). MR-Stream is a combined on- and offline technique with the cluster grid updated online and regular second stage macro-clustering of the grid combined with a tree pruning algorithm. In the extreme case of densely populated data space the offline components must visit each grid space,  $2^{DH}$  where  $D$  is the number of dimensions and  $H$  the granularity, or resolution, of the clusters.

A recently introduced technique, CODAS [18] demonstrated a new approach to clustering of continuous data streams into arbitrary shaped clusters. CODAS is a two stage technique with a micro-clustering first stage. The micro-clusters are designated as having an inner ‘kernel’ and outer ‘shell’ region. This removes the need for classifying micro-clusters as ‘edge’ or ‘non-edge’ and simplifies the micro-cluster joining calculations. As a result the technique has been demonstrated to be dimensionally stable with a time penalty in the order of  $(\frac{D}{100})$ , where  $D$  is the number of dimensions. Although CODAS is an online technique it does not allow for the clusters to evolve. That is to say that clusters, once formed, will remain. This means that in cases where the data stream evolves, i.e. forms different clusters at different times, CODAS does not update to remove the old micro-clusters. As a fully online technique, no data is stored and so it is not possible to use techniques such as windowing, or ageing of data to overcome this limitation. The technique presented in this paper builds on the underlying techniques of CODAS but stores the micro-cluster adjacency information in a graph structure to allow the macro-clusters to evolve and to provide additional information on the structure of the cluster-graph.

### 3. The Proposed Approach

Traditional offline clustering techniques for arbitrary shapes may categorize data samples as ‘core’ or ‘non-core’. However, this requires storage of the data samples and ever-increasing storage capacity which is prohibitive for online clustering. CEDAS stores only the information related to the micro-clusters and a graph structure recording the micro-cluster connections.

The following terminology is defined for the CEDAS approach:

1. Cluster Graph: the structure that defines which micro-clusters join to form which macro-clusters. This is stored by recording the intersects of each micro-cluster in ‘Edge’, together with the appropriate macro-cluster assignation in ‘Macro’.
2. Local density: the number of samples per micro-cluster
3. Macro-cluster: a cluster consisting of a number of intersecting micro-clusters.

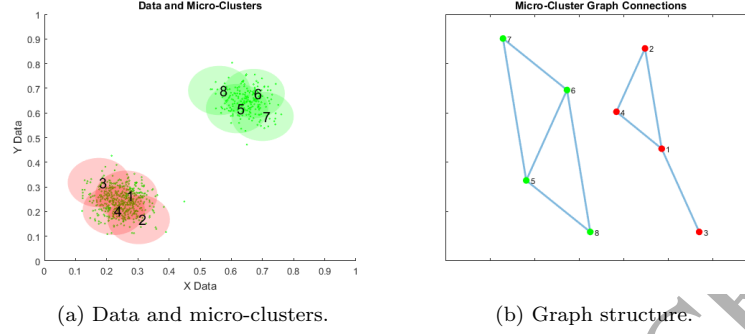


Figure 1: Example of the CEDAS algorithm micro-clusters and graph structure. The data together with two macro-clusters in red and green are shown in Figure 1a. Figure 1b shows the cluster graph structure with the nodes of the sub-graphs coloured according to the macro-clusters.

4. Micro-cluster: a micro-cluster with a local density above the threshold.
5. Outlier-micro-cluster: a micro cluster with local density below the threshold.
6. Sample: any data point in ' $D$ ' dimensions.
7. Threshold: the minimum number of samples within the micro-cluster radius of any sample to form a micro-cluster.

In general, CEDAS is a data-driven approach to divide the data space into kernel and shell regions based on a user defined radius,  $r_0$ . Each micro-cluster consists of a shell annulus region between radii  $\frac{r_0}{2}$  and  $r_0$  and a kernel region being  $r \leq \frac{r_0}{2}$ . Any micro-cluster above a given density threshold is considered for macro-cluster membership. Micro-clusters with kernel regions that intersect another micro-cluster shell region form macro-clusters. Micro-clusters above the threshold but with no intersections are also considered to be macro-clusters. Shell regions are considered to be edges of macro-clusters.

New data from the data stream will fall in to one of 3 regions:

1. empty space, where it will form a new, outlier-micro-cluster
2. a micro-cluster shell region, where it will be assigned to the cluster, the cluster count updated and the micro-cluster centre recursively updated to the mean of its samples.
3. a micro-cluster kernel region, where it will be assigned to the micro-cluster and the cluster count updated

The micro-cluster that has been modified, or created, by this process is then checked to see if the local density is above the threshold. If this is the case then this micro-cluster is checked for new intersections with other micro-clusters. If new intersections have been made then all the micro-clusters are linked and assigned to the same macro-cluster. This ensures that all linked micro-clusters



200 have the same macro-cluster reference and maintains arbitrarily shaped data space regions of macro-clusters in a fully online manner.

205 With this approach at any given time a data sample can be checked for its macro-cluster membership, any new sample is immediately clustered and outliers are identified as members of outlier micro-clusters. It is the graph structure for storing the micro-cluster intersections that forms the basis of the advance from CODAS to CEDAS as this allows rapid merging and division of macro-clusters.

### 3.1. CEDAS Algorithm Description

210 There are 4 distinct steps for the full algorithm including initialisation and a *Step 0* is included where the user determines the parameters for the algorithm:

0. Parameter Selection
1. Initialization
2. Update Micro-Clusters
3. Kill Clusters
- 215 4. Update Cluster Graph

220 A description of each of the key algorithm steps is provided here. The algorithm runs these sections sequentially for each data sample, if required. Not every section is required each time and the conditions for running each section are described. A Matlab implementation of the algorithm is available on Github [17].

#### 3.1.1. Parameter Selection

CEDAS requires a number of parameters to function, *Decay*, *radius* and *Minimum Density Threshold*. The values for these parameters are application dependent and suitable values can be selected as follows:

- 225 1. *Decay*: the decay is specified as a number of data samples in the implementation used in the preparation of this paper. If the data examined has a regular sample rate this is directly related to the length of time over which the data is to be examined, e.g. at a sample rate of  $1Hz$ , to examine the data over a 28 day period the Decay would be 2,419,200, at  $0.1Hz$  for 230 7 days, 60,480. In the case of an irregular data rate a time based decay value could be used.
- 235 2. *radius*: the radius of the micro-cluster is selected based on expert knowledge of the application. With any set of data there are distances between data samples. There is a maximum distance between data beyond which an expert will consider that the data belongs to a different cluster and this value is the maximum allowable radius, i.e. the radius should be set to the minimum allowable gap between macro-clusters. Using a radius below this value has little effect on the overall macro-cluster beyond compiling them from a greater number of micro-clusters and smoothing the profile of the macro-cluster. There is an effective lower limit to the radius below 240 which it will not contain enough data samples for a micro-cluster to form.

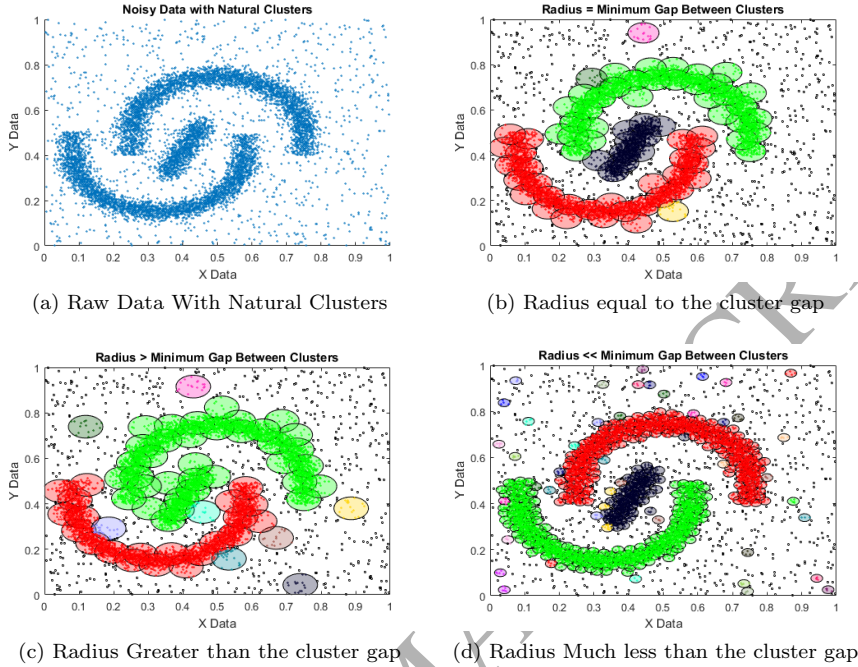


Figure 2: Demonstration of varying CEDAS radius selection. Figure 2a shows raw data with noise and natural clusters. Figure 2b shows the cluster results with radius equal to the minimum gap between the clusters, Figure 2c shows the results of having a larger radius than the minimum gap and Figure 2d shows the effect of a much smaller radius. Thus radius is set by the user and should be less than the maximum dissimilarity data can have and still be considered a part of the same cluster.

Visual examples are shown in Figure 2 where figure 2b shows successful clustering with the radius equal to the minimum gap between natural clusters, Figure 2c shows how increasing this value determines that two of the clusters are not different enough to be considered separate and they become merged. Figure 2d shows how reducing the radius to half of the minimum gap has little effect on the macro-cluster results.

3. *MinimumDensityThreshold* is required to differentiate clusters from background noise. The value should be set based on expert knowledge as to the level of data required to be considered valid, natural clusters.

### 3.1.2. Initialization

This creates a structure to store the information related to each micro-cluster and takes place with the first data sample. The '*Centre*' defines the location of the micro-cluster in data space. '*Count*' stores the total number of data samples that have been allocated to the micro-cluster. The value of '*Count*' is recorded to allow recursive updates to the micro-cluster centre. '*Macro*' is a

---

**Algorithm 1:** CEDAS: Initialization

---

**Input:**  $x, r_0$

Create micro-cluster structure containing:

$C_1(Centre) = x$

$C_1(Count) = 1$

$C_1(Macro) = 1$

$C_1(Energy) = 1$

$C_1(Edge) = 1$

Set number of micro-clusters to 1

Set modified micro-cluster number, for use updating the graph structure.

---

reference to the macro-cluster to which this micro-cluster belongs. The value of ‘*Macro*’ is the same for all micro-clusters in the ‘*Edge*’ list. ‘*Energy*’ is a value used to determine the length of time since a micro-cluster received new data. The decay algorithm reduces this value and is discussed later. ‘*Edge*’ is a list of intersecting micro-clusters, if a micro-cluster has no Edge list then it is a macro-cluster by itself. In graph theory terminology the micro-cluster number paired with each intersect constitutes an ‘edge’ of the form  $\{mC_c, mC_i\}$ , where the first term is the current micro-cluster and the second term is the intersecting micro-cluster.

### 3.1.3. Update Cluster

This part of the algorithm updates the micro-clusters when a new data sample arrives.

---

**Algorithm 2:** CEDAS: Update Micro-Cluster

---

**Input:**  $x, C, r_0$

find distance to nearest micro-cluster centre,  $d_{min}$

**if**  $d_{min} < r_0$  **then**

    reset micro-cluster *Energy* to 1

    increment number of samples contained in micro-cluster

**if** *data is within micro-cluster shell* **then**

        | recursively update micro-cluster centre

**end**

**else**

    | Create new micro-cluster

**end**

---

The algorithm checks whether the new data sample belongs to any current micro-cluster. If it does not then a new micro-cluster is created. If the data sample is within a current micro-cluster then the *Energy* of that micro-cluster is re-set to 1 and the number of data samples it contains is incremented. A further check is made to find if the sample lies within the kernel or shell of the micro-cluster. If it is in the shell region then the centre of the micro-cluster is

275 recursively update to the mean of the data samples in the shell. Only updating  
the centre if the data lies within the shell has the effect of prevent a single  
micro-cluster endlessly following drifting data by limiting its movement. (The  
same effect is also achieved by only updating the centre if the sample lies in the  
kernel.)

#### 280 3.1.4. Kill Micro-Cluster

This part of the algorithm reduces the energy of the micro-clusters and  
removes them if the energy has fallen below zero.

---

**Algorithm 3:** CEDAS: Kill Micro-Cluster

---

**Input:**  $C, Decay$   
Reduce all  $C(Energy)$  by  $Decay$   
**if** Any  $C(Energy) < 0$  **then**  
    Remove micro-cluster  
    Remove all edges containing the micro-cluster  
    Decrement the number of micro-clusters  
**end**

---

First, all the micro-cluster energies are reduced by the decay amount. Then,  
if any micro-cluster energies are below zero, they are removed. All edges that  
285 refer to this micro-cluster are also removed and the total number of micro-  
clusters is reduced.

#### 3.1.5. Update Micro-Cluster Graph

---

**Algorithm 4:** CEDAS: Update Graph

---

**if** A micro-cluster has been modified **then**  
    **if** the micro-cluster edge list has changed **then**  
        Set a new macro-cluster number throughout the graph  
    **end**  
**end**  
**if** Any micro-clusters have died **then**  
    Set new macro-numbers for the sub-graphs of its previous edges  
**end**

---

This section only makes any changes if either:

1. a new cluster has been created and reached the minimum density threshold
2. a cluster centre location has been modified
- 290 3. a cluster has died and been removed.

First the changes are made to any micro-cluster that has been modified by  
either having its centre location moved or by virtue of being a micro-cluster  
that has newly reached the threshold. In either case the graph edges for that

micro-cluster may have changed. If the edge list has changed then the new graph has its macro-cluster number set to a new value.

The changes made by any micro-cluster that have died out are then addressed. Any micro-clusters that the dead micro-cluster used to have an edge with have their graphs updated with a new macro-cluster number.

#### 4. Experimental Results and Discussions

This section analyses the performance of the CEDAS algorithm and presents the results and discussion across a range of experiments. In subsection 4.1 we validate the ability of CEDAS to accurately deal with data drift, cluster separation, cluster merging and noise over time. We then compare the speed and accuracy with alternative techniques CluStream, DenStream and MR-Stream across high dimensionality data in subsection 4.2. In subsection 4.3 we compare CEDAS to these techniques with regard to complexity, processing speed, cluster quality and memory efficiency. Finally in subsection 4.4 we apply the CEDAS algorithm to a real data stream from the London Air Quality monitoring system to demonstrate how evolving clustering can aid data mining of data streams containing short term drift, long term drift and short and long term anomalies.

##### 4.1. CEDAS Functionality with Cluster Separation, Cluster Merging, Drift and Noise

A 3D data stream consisting of 2 Mackey-Glass time series is presented as a data stream. The data stream is a pair of solutions of the Mackey-Glass non-linear time delay differential equation [21, 13]. shown in equation 1.

$$\frac{dx(t)}{dt} = \frac{ax(t-\tau)}{1+x(t-\tau)^{10}-bx(t)} \quad (1)$$

The equation is solved numerically at discrete time steps using the 4th order Runge-Kutta method using different values for  $a$  and  $b$  to create  $x$  and  $y$  values as shown in Figure 3a. For each time step 10 random data samples were created around the core value to provide a data stream of 40,020 samples illustrated in Figure 3b. Early in the data the values of both data streams are coincident. They later separate and come together at various times. We would expect that ‘recent’ data will produce a changing number of macro-clusters, as the data separate and rejoin, and that an online, evolving clustering technique will detect these changes as they occur. A further data set was created by adding additional data of random noise samples at every 5th time step creating a dataset of 44,022 samples. These data are used to test the robustness of CEDAS to detecting the clusters in a noisy environment. By presenting the data sequentially we create a continually evolving data stream rather than a data stream of similar values with sporadic variation, such as the KDCCup data set below. This tests the ability of the algorithm to add, merge and separate macro-clusters in a continuously evolving environment.

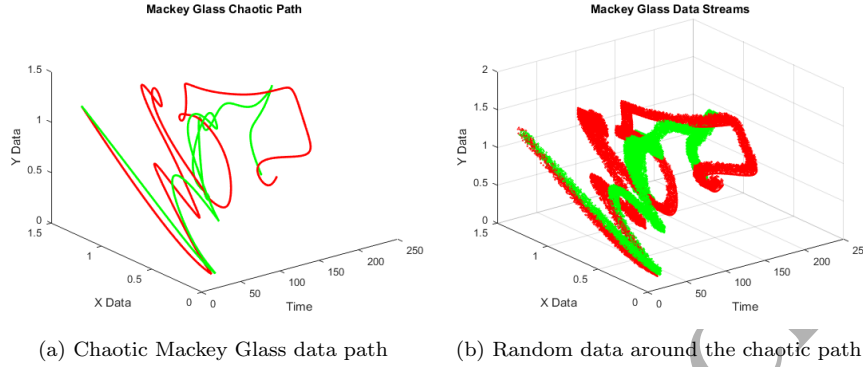


Figure 3: Illustration of the Mackey-Glass data sets showing a) the chaotic path b) the data stream created around that path. The two Mackey-Glass streams are shown in red and green. When considering the data over the previous ' $N$ ' samples the data may form separate streams, two clusters, or streams that are joined at some point, a single cluster.

To validate the correct functionality of CEDAS the algorithm was applied to the Mackey-Glass data streams using  $Decay = 1,000$  samples,  $Radius = 0.05$  and  $MinimumThreshold = 15$ . The decay is a suitable time period for investigation such that the macro-clusters will be large enough to visualise, and the two Mackey-Glass data streams will be merged and separated for sufficient time to indicate the correct operation of CEDAS. The radius is selected such that the 'width' of the data streams is encapsulated minimising the plotting time for multiple micro-cluster spheres. The minimum cluster size is such that it is larger than the expected density of the noisy data ensuring that the noise remains as outliers.

The data was presented to the CEDAS algorithm 1 sample at a time to imitate an online data stream and the results plotted at each time step to create a video of the results. The CEDAS algorithm was used to detect and report in the plot title the following information:

1. Definite Clusters: these were defined as clusters containing  $> 15$  data samples and  $> 1$  micro-cluster. These are settings specific to the investigation to interpret the algorithm results and are not algorithm parameters.
2. Outlier Clusters: these were defined as containing  $> 15$  data samples all contained in 1 micro-cluster. These are also specific to interpretation of the results and not algorithm parameters.
3. Last Change: the time period at which the last change in the number of Definite Clusters occurred. This information was recorded to allow the state at that time to be reproduced.

#### 4.1.1. Cluster Separation and Merging

Using the clean Mackey-Glass data stream the sample number at which a change in the number of macro-clusters was detected was stored. After the

Table 1: Comparison of trigger points with and without noise. The trigger points in brackets with noise are short term and caused by the effect of noise in moving the micro-cluster positions briefly.

Trigger Points		
Group	Without Noise	With Noise
1	a	a
1	b	b (c)
3	c	d (e, f, g, h)
4	d	i
5	e	j
6	f	k
7	g	l
8	h	m
9	i	n
10	j	(o) p

analysis data was plotted with data the data coloured differently each time the number of macro clusters changed. This is shown in Figure 4a.

After the initial settling period (red), it is seen that at each colour change the number of clusters in the data contained in the preceding decay period has changed. For example, in the green period the data was contained in a single cluster. At the time the colour changes to black, the data in the previous 1,000 samples had just separated to 2 separate macro-clusters. When the colour changes to magenta, now the previous 1,000 samples create 1 macro-cluster. The colours of the data do not represent the clusters themselves, but represent the period preceding the time at which the number of macro-clusters changes.

#### 4.1.2. The Effects of Noise

To test the effects of noise on CEDAS the Mackey-Glass dataset is used with a random noise sample added every 5 data samples as described above. The random nature of the noise will have some effect on the initial positions of micro-clusters if the noise falls within them. This increases the likelihood of an initial micro-cluster separating from the main macro-cluster group. If this occurs then the number of macro-clusters may change briefly. This would give the appearance of false positives when compared with the results from dataset without noise. These additional clusters are in fact present at that time and it is accepted that the noise has in fact changed the clustering.

The results are shown in Figures 4a and 4b. Figure 4b illustrates that *c, e, f, g, h, o* are triggered by the noise and could easily be discounted based on the number of samples, if required. The trigger points without noise can be matched to those with noise as shown in table 1. These are discussed in the following sub-sections.

#### 4.1.3. False Positives

With any online technique, apparent changes at some point in time may turn out to be irrelevant at a later time. An example of such soon-to-be-irrelevant data anomalies are those that result from the added noise. Rather than calling

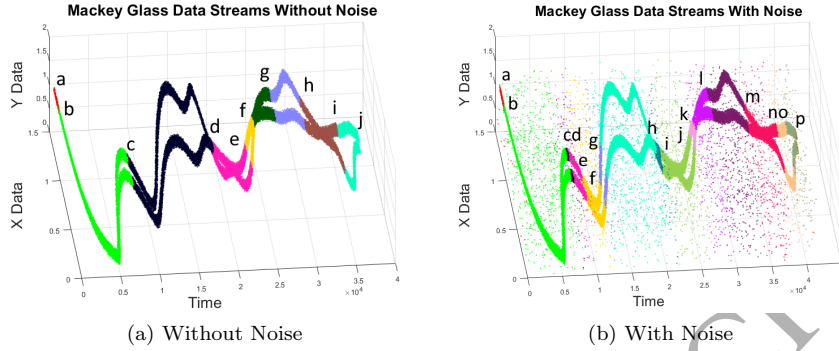


Figure 4: CEDAS Auto Change Detection, changes in colour represent changes in the number of clusters. Thus in figure 4a while the data is coloured green previous ' $N$ ' samples form a single cluster, joined at the beginning. At the point the data colour changes to black, the data in the previous ' $N$ ' samples has separated into two clusters. It should be noted that the colours of the data are not the clusters themselves, but represent the time periods during which the data forms different numbers of clusters. The changes detected without noise are also detected with noise with the additional changes caused by temporary separate micro-clusters before they rejoin the main clusters.

these 'false positives', they could be considered as 'temporary or short-term true positives'. In the event these are caused by temporary misplacement of micro-clusters caused by noise, which are rapidly re-absorbed into the macro-cluster, then these addition clusters will have an unusually short lifespan, i.e. considerably shorter than the set decay period. In this way any triggers that are within a user-defined short time span from a previous trigger could be discounted if required. However, this is not always desirable, as even short term anomalies may be of interest. They may, for example, indicate the start a general drift in the data.

#### 4.1.4. False Negatives

With appropriate settings for decay time and micro-cluster radius, false negatives do not occur. It must be remembered that a different decay time will create different times for cluster separation. This is not indicative of false negatives, but rather a deliberate function of the technique to consider clusters based on data within a defined time frame.

#### 4.1.5. True Positives

As we have demonstrated here, all changes to clusters are correctly detected. With the noisy dataset, although we have some temporary true positives, as discussed above, CEDAS has successfully detected the same true positives as with the clean dataset as shown in table 1.



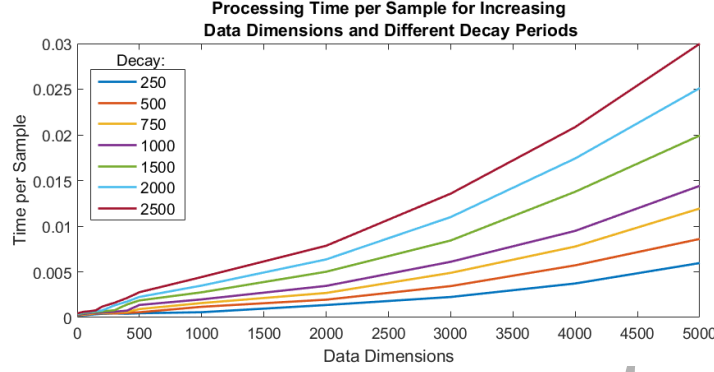


Figure 5: Plot of mean processing time per sample in seconds for varying data dimensionality. Each line represents the processing time for different decay periods which create a proportional increase in micro-clusters, e.g. the top, red line represents the processing time per sample for a decay period of 2,500 samples for data with dimensions from 1 to 5,000.

#### 4.1.6. True Negatives

If we consider the definition of a ‘true negative’ to be that ‘no changes in macro-clusters are detected when there are none’ then this occurs with every sample that does not create new clusters.

#### 4.2. High Dimensional Data Test

The dataset used in this section comprises three helical data streams, two of which join mid-way through the test while the other stays separate. These data streams are moved through a range of multiple dimensions to examine the time variance of the analysis with higher dimensional data. The data was analysed using CEDAS with a range of values for *Decay* and settings of *InitialRadius* = 0.05 and *MinimumThreshold* = 4. As the aim of this experiment is to test the speed penalty across high dimensional data and not to test the efficacy, a-priori knowledge of the data streams was used to ensure valid clustering occurs. *Decay* was set at a reasonable number of samples to ensure macro-clusters of a suitable size to demonstrate the effectiveness of the technique. The radius is set smaller than the width of the helices to ensure multiple micro-clusters at all times, and below the minimum expected gap between natural clusters. The minimum threshold was set to 4 to restrict micro-clusters from forming on the very edge of the natural clusters. The data set is then moved into higher dimensional data space by adding additional dimensional data coordinates. By projecting the data back into 3 dimensions the clustered data can be plotted and the results of cluster membership checked while increasing the complexity of the clustering calculations.

##### 4.2.1. Speed and Dimensionality Comparison

By utilising hyper-spheres for micro-clusters the cluster joining technique checking for micro-cluster overlap is much simpler than, e.g. hyper-ellipsoidal

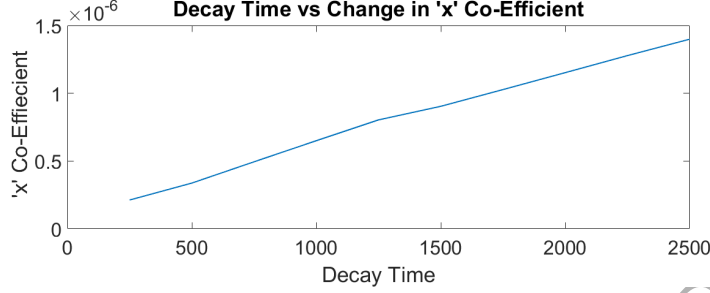


Figure 6: Comparison of the processing time per sample with the decay time showing a linear relationship between processing speed and decay time. For the data used in this example the decay time is directly proportional to the number of micro-clusters. Where a longer decay time does not results in additional micro-clusters, then the time per sample remains constant. In practice the processing time will lie somewhere between the two.

micro-clusters. Micro-clusters are joined if the edge of the core hypersphere  
 intersects another hyper-sphere shell. This requires only a comparison between  
 the euclidean distance between cluster centres and the sum of the micro-cluster  
 radii. Therefore, the only calculation that is dimensionally dependent is the  
 euclidean distance with complexity  $O(D)$  where ' $D$ ' is the number of dimensions.  
 The relationship between the number of data dimensions and processing time  
 per sample is linear.

With each new data sample being assigned to a single micro-cluster it is  
 only necessary to check the intersections for that micro-cluster and then only  
 if the micro-cluster centre has been modified, or a new micro-cluster has been  
 created. This further reduces the required number of calculations. The radii  
 of the micro-clusters is constant and so we only need to compare the euclidean  
 distance between the changed micro-cluster and all others with  $1.5r_0$ .

The relationship between the number of data dimensions, decay period and  
 calculation time is plotted in Figure 5. Using Matlab's curve fitting toolbox the  
 relationship between the data dimensions and time per sample was tested for  
 different decay times. In the case of an evolving data stream with continuous  
 drift the decay time is also proportional to the number of micro-clusters. To  
 investigate the relationship between decay time, and so the number of micro-  
 clusters, and run time the coefficients of the ' $x$ ' term in the linear equations  
 are plotted in Figure 6 and show an approximately linear relationship. These  
 results concur with the predicted linear time penalty for both the number of  
 dimensions and the number of micro-clusters.

By comparison, Figure 7 shows the relationship between processing time and  
 dimensionality with the same data set for comparison with both DenStream [7]  
 and Clustream [1]. We used the Massive Online Analysis [5] implementation  
 running on R3.2.2 in RStudio 0.98.1102 analysing the same helical high dimen-  
 sional dataset as for CEDAS. CluStream was also limited to a maximum of 100  
 micro-clusters. For both of these techniques, two tests were run using a decay

time of 1,000 samples:

1. Both DenStream and CluStream without carrying out the 2nd stage re-clustering until the end of the data stream.
2. We approximated a fully online technique by carrying out the 2nd stage clustering technique at frequent intervals - every 100 samples for DenStream and every 10 samples for CluStream.

For the DenStream 2nd stage re-clustering DBScan [12] was used as implemented in the ‘R’ package by Hahsler [14] to allow for arbitrary shaped macro-clusters to form in a similar manner to CEDAS. The results shown in Figures 7a and 7b are for test 1 and the results shown in Figure 7c and 7d are for test 2. Without 2nd stage re-clustering both DenStream and CluStream are faster than CEDAS for low dimensionality data. The break even point is approximately  $12D$  for CluStream and  $220D$  for DenStream. When the second stage re-clustering of the micro-clusters is done frequently enough to approximate fully online analysis there is significant time penalty for both DenStream and CluStream. In both cases CEDAS is noticeably faster than both DenStream and CluStream and suffers significantly less time penalty for increasing data dimensionality.

#### 4.3. CEDAS Speed, Purity, Accuracy and Cluster Validation

To further test the CEDAS algorithm in a different environment the KDD-Cup99 [16] dataset was used as a data stream by presenting the data to the algorithm sequentially. The data set consists of approximately 5 million samples in the full data set, 500,000 samples in the 10% reduced set, simulating network intrusion attacks on a military installation. The dataset has 42 features and information to classify the data into 22 attack types in addition to the normal network traffic. This data is used to determine the cluster purity and memory use for comparison with alternative techniques and also to validate the clustering results in relation to the number of attack types which occur in a time period.

##### 4.3.1. Speed and Cluster Quality

The KDDCup99 data stream is a popular dataset for testing evolving clustering algorithms such as eClass [2] and it is used here to allow direct comparisons with D-Stream and MR-Stream purity results presented by Wan et al [28]. Two sets of results are presented. The first is the same analysis used by Wan et al. of creating 500 time intervals spaced at 1K samples and placing these into groups of 25 and taking the mean cluster purity over these groups of 25. Taking the mean of a set of results can disguise individual poor results and so the cluster purity for CEDAS at each of the 500 time intervals is also provided. These results are shown in Figure 8.

It should be noted that the mean cluster purity alone, as defined by equation 2, may be a poor measure by itself.

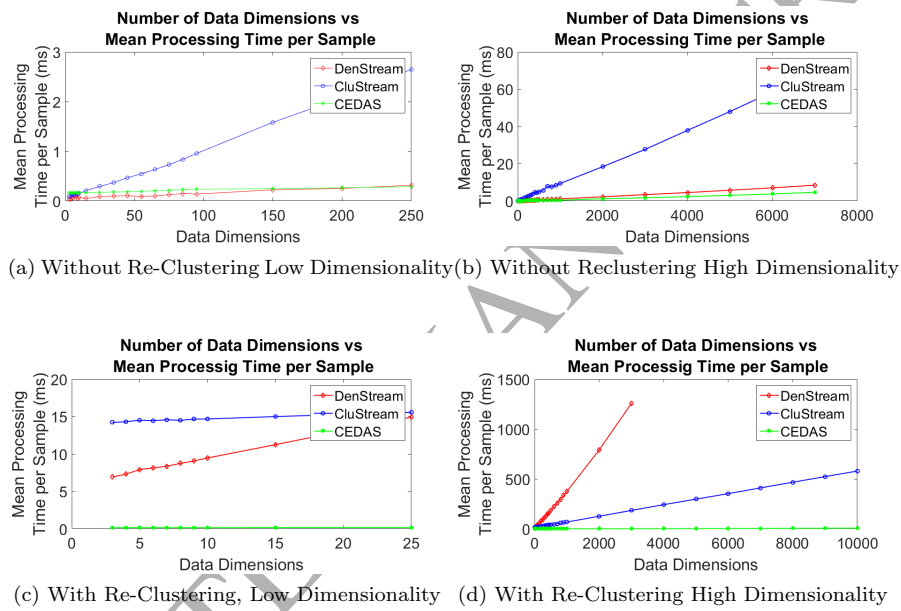


Figure 7: Typical analysis time per sample for DenStream, CluStream and CEDAS across various dimensional data. a) and b) show CluStream and DenStream without 2nd stage re-clustering until the end of the data stream. c) and d) show DenStream and CluStream with frequent 2nd stage re-clustering. In all plots CEDAS is shown in green. DenStream and CluStream have a faster 1st stage clustering, but for fully online clustering CEDAS is shown to be faster.

$$mean\ purity = \frac{\sum_{i=1}^N \frac{|C_i^d|}{|C_i|}}{N} \times 100\% \quad (2)$$

$$accuracy = \frac{\sum_{i=1}^N |C_i^d|}{\sum_{i=1}^N |C_i|} \times 100\% \quad (3)$$

Here  $C_i$  is the number of samples in a cluster,  $C_i^D$  is the number of these samples assigned to the dominant class and  $N$  is the number of clusters. In cases where a high number of samples are contained in one cluster with low purity, yet few samples are contained in a high number of clusters with high purity the result is a high mean purity even though most samples are incorrectly assigned. Equally, the reverse is true when few clusters are present, if 99% of the data is correctly assigned in one cluster and two sample are contained in a second, one of which is mis-assigned the mean purity looks poor. In Wan et al. the relevance of this measure is further reduced by taking the mean of these means and so the purity measure is included here for comparison to Wan et al. only and not to attach any particular significance to the result. The cluster accuracy measure as defined in equation 3 is presented in Figure 8d which is a measure of the number of clustered samples that have been correctly assigned to the dominant class. By using both the purity and accuracy measures the quality of the clustering can be stated with greater confidence.

The results of the quality analysis are shown in Figure 7. Although the purity at time period 145 is 73%, the mean over the 25 time periods this is 96%. Using the two time periods selected by Wan et al, 27 and 52, we see that the CEDAS purity was 96% and 99.85% compared with MR-Stream at 97.5% and 92% respectively. It is interesting to note that at time periods 26 and 28 CEDAS purity is 100% suggesting that CEDAS adapts quickly to this variation. Using the 25 time periods measure favoured by Wan et al. we see that CEDAS mean purity exceeds that of MR-Stream. When considering the accuracy measure we note that at the time periods 27 and 52 the accuracy measurements are 98.5% and 99.98% respectively. This indicates that nearly all the samples are correctly assigned to the dominant clusters, but the purity is reduced due to few incorrectly assigned samples in clusters with few members. The accuracy of CEDAS remains close to 100% at all times except for 3 single occasions where it drops to around 90% and 2 at around 95%.

Figure 8a shows the results provide by Wan et al. for the mean purity over 25 time steps for both D-Stream and MR-Stream.

Having established via the purity and accuracy measures that the clusters are meaningful it is useful to see if they demonstrate any results of interest. To do this the number of clusters in a time period are compared with the number of classes given in the data. The plot of these is given in Figure 9 where it can be seen that each time there is a rise in the number of classes, i.e. attacks, the number of clusters also rises. Given that these clusters have high purity, and the accuracy of clustering is also high, these additional clusters must contain attack vectors unique to each type of attack. There are 50 time periods with attack

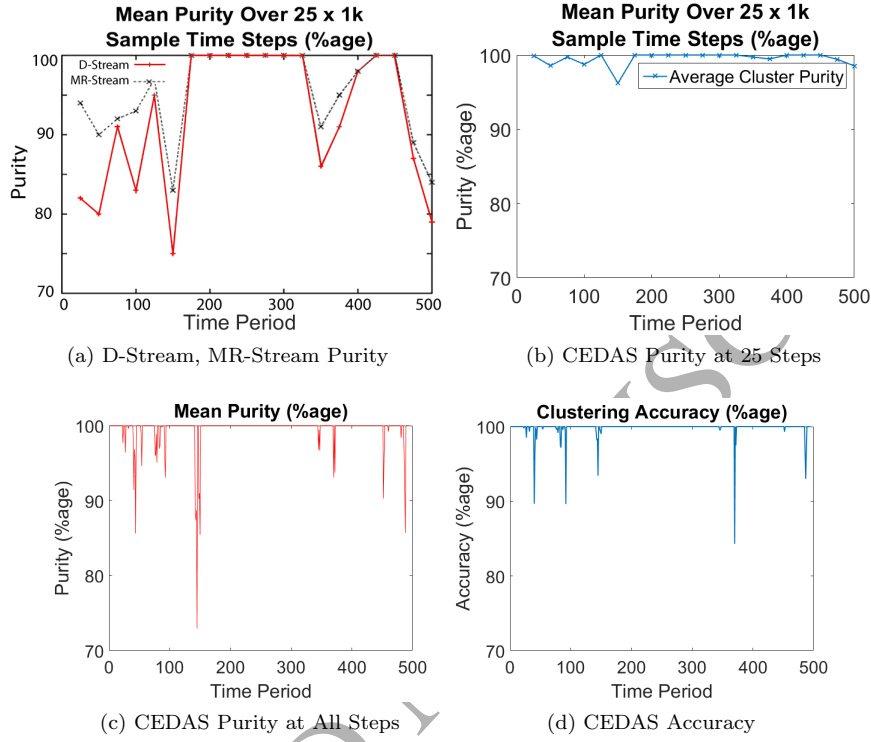


Figure 8: (a) Plot of mean cluster purity (from [28]), (b) Mean cluster purity for CEDAS by the same measure as Wan et al. [28]. (c) Cluster purity at each time step showing instances of reduced mean purity. (d) CEDAS accuracy measure.

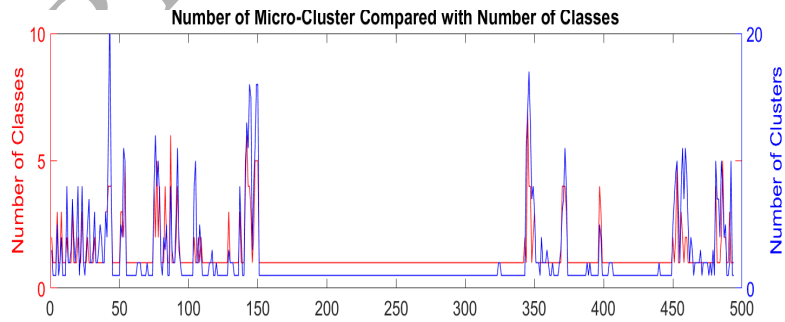


Figure 9: Plot of the number of classes of attack and the number of clusters found by CEDAS in each time period. The number of clusters is proportional to the number of classes throughout.

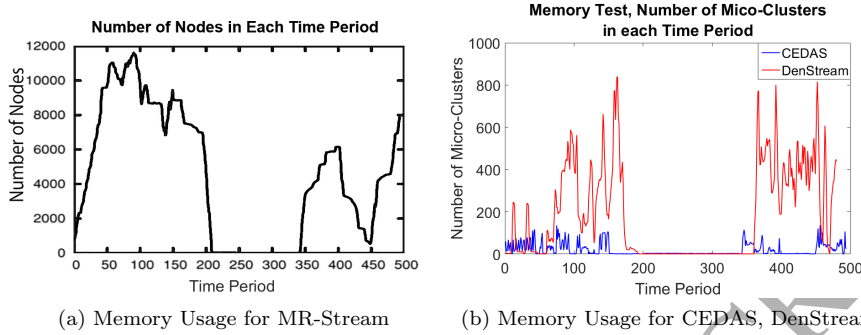


Figure 10: Plot of the number of nodes or micro-clusters, which equates to memory use, for MR-Stream (from [28]), DenStream and CEDAS. CluStream is not shown as it uses a maximum number of micro-clusters set by the user. CEDAS shows the lowest memory use.

vectors present and these are detected 100%. As discussed above, occasional separated micro-clusters are a feature of evolving techniques and providing they are short-lived and re-absorbed into the main clusters they can be ignored with reasonable confidence. When the number grows beyond 1 sample per cluster, however, they may be indicative of possible attacks. Thus with a threshold of 1 we have 20 false positives. However increasing the threshold to 2 to allow for occasional separated micro-clusters reduces this Figure to 4, and a threshold of 3 reduces this to a single instance. This compares favourably with a mean number of clusters per attack of 8.2.

#### 4.3.2. Memory Efficiency

To demonstrate the efficient memory use of CEDAS, we compare the storage required by MR-Stream and DenStream with that required by CEDAS when clustering the KDDCup99 datastream. The results presented by Wan et al. for MR-Stream are shown in Figure 10a and, when the data stream is evolving and has variety, we see that MR-Stream reaches Figures in the thousands of nodes with a peak approaching 12,000. By contrast, the number of micro-clusters required by DenStream and CEDAS for the same data stream are shown in Figure 10b. DenStream has a mean value of 181 and maximum of 839 whereas CEDAS has a mean of 20 and peaks at 137. This demonstrates the significant memory saving of micro-clusters over grid based techniques. Even allowing for the CEDAS cluster description consisting of 5 values there is significant saving over MR-Stream.

#### 4.4. Data Mining of Atmospheric Data Streams Using CEDAS

In this section we apply the technique to data from the Kings College London Air Quality Website [11]. The data is from one monitoring site, Westminster Marylebone, and we use 2 dimensions, labelled  $NO_x$ ,  $PM_{10}$ . Here and throughout,  $NO_x$  is defined as the reactive oxides of Nitrogen, primarily  $NO$  and  $NO_2$ ,

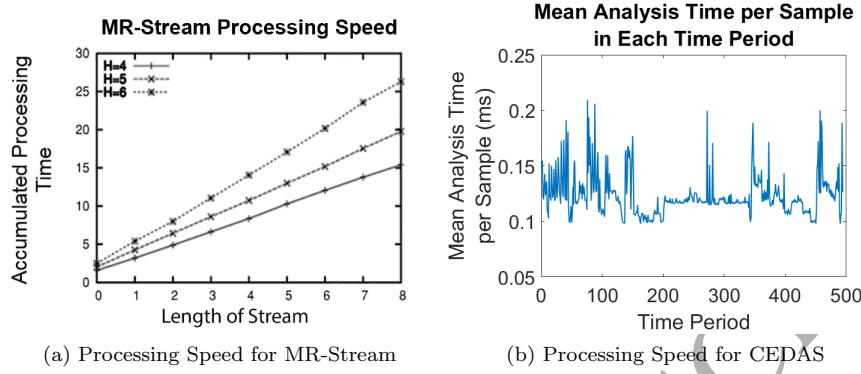


Figure 11: Plots of the processing speed of 10a MR-Stream accumulated time (from [28]) and 10b CEDAS mean time per sample.

and  $PM_{10}$  is defined as the mass concentration of microscopic airborne particles with aerodynamic diameter of  $10\mu m$  or above. The data, which is recorded operationally to monitor breaches of air pollution legislation [25] and to inform the public of adverse air pollution conditions, is captured at 15 minute intervals and ranges from 1<sup>st</sup> January 2010 to 30<sup>th</sup> December 2014 for a total of 87,600 samples. This data is used to test CEDAS ability to differentiate short and long term anomalies and follow the temporal drift of real data.

To allow for clustering to take place the data is normalised to a suitable range relative to the micro-cluster radius,  $r_0$ . Here the range was based on the data available in the dataset and scaled to 0 – 1. The data had an actual range from  $min = 7.200$  to  $max = 1,447, ppbv$  (parts-per-billion by volume) for  $NO_x$  and  $min = -0.9$ ,  $max = 422.8 (\mu gm^{-3})$  for  $PM_{10}$  and so predicted ranges of 0 to 1500 and 0 to 200 respectively were used. The scaling introduced by this normalisation has an effect on the local density, joining and separation of micro-clusters and so expert knowledge is required to find suitable values for scientific research involving the cluster results.

We describe how anomaly detection differs between long-, medium and short-term analysis and how CEDAS copes with such variation. To demonstrate this we define ‘Short Term’ as being 7 days and ‘Medium Term’ as being 28 days and ‘Long Term’ as being one year. The *decay* values used correspond to the number of data samples collected in the respective time period. For the *radius* value we have used a value of 0.05. This value is arrived at by looking at historical data and estimating the distance from the main natural cluster to data that we would consider an outlier. The definition of what is considered different enough to be an outlier is at the discretion of the expert user. The analysis carried out here is robust to a range of *radius* values with little change in the macro-clusters or their visual appearance. We are interested to see all data space regions containing data, including single outliers and so we set the minimum threshold to 1.



We use data collected between 2010 and 2014 inclusive. The data was presented to the CEDAS algorithm sequentially, in  $NO_x$ ,  $PM_{10}$  pairs, to mimic an online data stream. The micro-clusters were plotted and the transparency of the micro-clusters set according to the value of the Energy in each. In this way if anomalous data appears for a short period of time the cluster adjusts, but it fades over the subsequent time period providing an online visualisation of the Energy of the micro-clusters. This provides a clear visual indication of CEDAS adapting to the changes in the data stream and following long term and short term drift. By using different decay times we demonstrate the different clusters that are created and discuss how this can be useful to investigate different time periods for drift and shift.

In Subsection 4.4.1 we consider how the use of a short decay period can reveal short term data drift that would be disguised in medium term decay periods. Subsection 4.4.2 describes the use of medium term decay periods to investigate possible seasonal variations. Finally, in Subsection 4.4.3 demonstrates how medium term decay periods can be used to investigate long term variations. In this paper we only provide visual indications of how CEDAS reacts to the evolving data stream. Potential numerical analysis of the clustering results is discussed in Section 5 Conclusions.

Videos of the CEDAS cluster analysis can be found in the supplementary material.

#### 4.4.1. Short Term Drift and Anomalies

Using a decay period equivalent to 7 days of data we can detect the changes in  $NO_x$  and  $PM_{10}$  over time. Sample plots are shown in Figure 12 (a)-(c) showing the cluster analysis at 3 different dates for the preceding 7 days. The data for the preceding 28 day period, for the same dates, is shown in Figures 12 (d)-(f).

We can see that the 7 day period preceding 24/03/11 is markedly different from the 7 day period preceding 06/02/11. Despite these difference in the 7 day data, by comparing the plots (d)-(f) we can see that overall, for the preceding 28 day periods the spread of data values has been more consistent. The data shown in the black and green clusters of the 7 day analysis in 12 (b) may be considered anomalous for that week, but in Figure 12 (e) we see that it is not unusual over the preceding 28 day period. However, data such as that in the yellow and magenta cluster of 12 (b) is seen to still be anomalous over the 28 day period, Figure 12 (e), where the clusters are now coloured khaki and blue.

This demonstrates that, by selecting suitable decay periods, the clustering results from the proposed algorithm provides relevant analysis of how data behaves over different time periods and how CEDAS can follow these changes in a fully online manner.

#### 4.4.2. Medium Term Drift and Anomalies

The plots in Figure 13 are the cluster results for a 28 day decay period taken at different dates throughout the year. Over the 5 year period of the data streams this approximate pattern is repeated each year. The primary variation

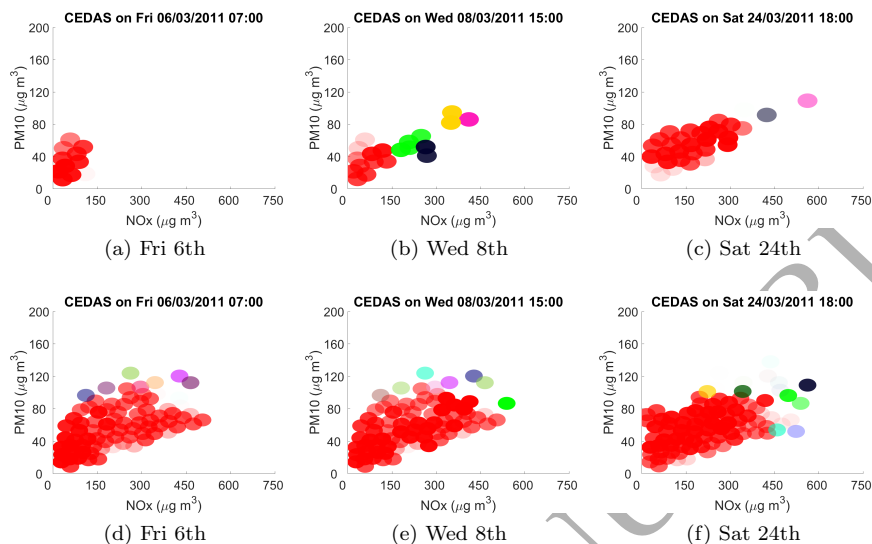


Figure 12: Sample plots of short term decay periods (a)-(c) and medium term decay periods (d)-(f). The short term variations indicated in (a)-(c) show the data varies over different 7 day periods. The medium term variations in (d)-(f) show that the data over the 28 day periods is more consistent and disguises the 7-day variation.

is not in the maximum, minimum or range of either  $NO_x$  or  $PM_{10}$  but rather in the range of the  $PM_{10} : NO_x$  ratio. This is particularly noticeable when comparing, e.g. March and July where at any given value of  $NO_x$  the range of  $PM_{10}$  values is greater in March. Anomalous data can also be seen in March indicating that some unusual events are present.

This demonstrates the ability of CEDAS to follow such seasonal drifts, if they exist, and find data that is anomalous within that local time frame.

#### 4.4.3. Long Term Drift and Anomalies

For long term changes, i.e. changes across years, the data could be analysed in multiple ways. For example, the data could be clustered on the full 365 day decay period. However, as we have already indicated in the Subsection 4.4.2 there are variations within that year which may be hidden in the way described in Subsection 4.4.1. With this information it is reasonable to consider an analysis of 28 day decay periods, at the same date, for subsequent years. Examples of these cluster results are provided in Figure 14 and shows the results for data of the 28 days preceding 01/04 for the years 2010-2015.

The shape of the main cluster can be seen to vary between years indicating the changes in data values. Anomalies are indicated and are for the particular month and year under consideration. In all cases we see some relatively minor anomalies with values that are slightly different from the main cluster. These could be symptomatic of the data undergoing normal drift and changes. March

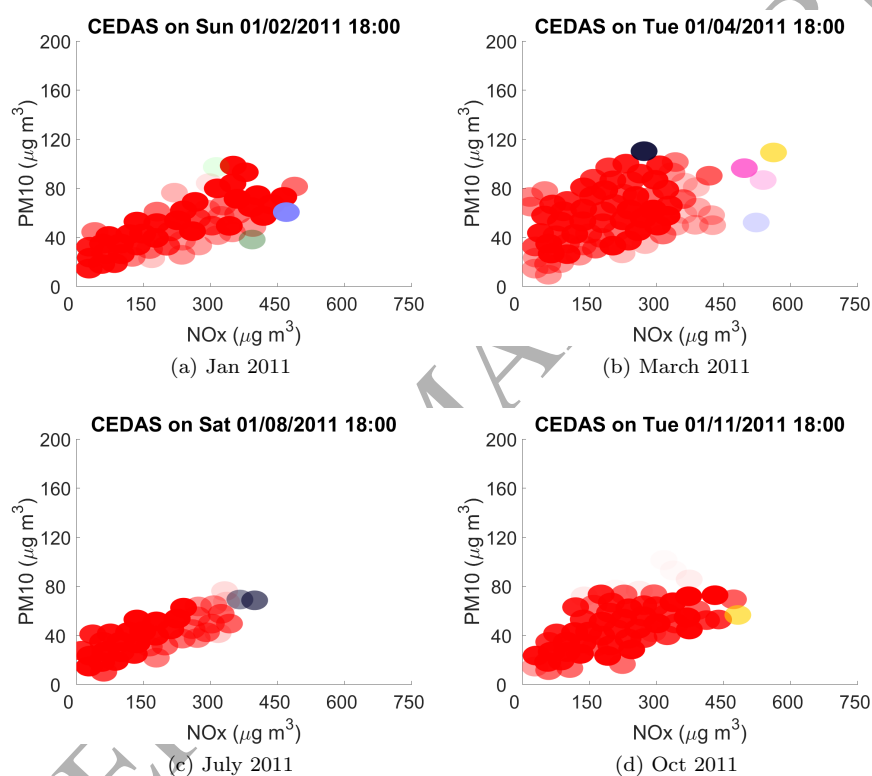


Figure 13: Plots of CEDAS clustering for a 28 day decay period showing a variation of the data spread at different dates during a single year.

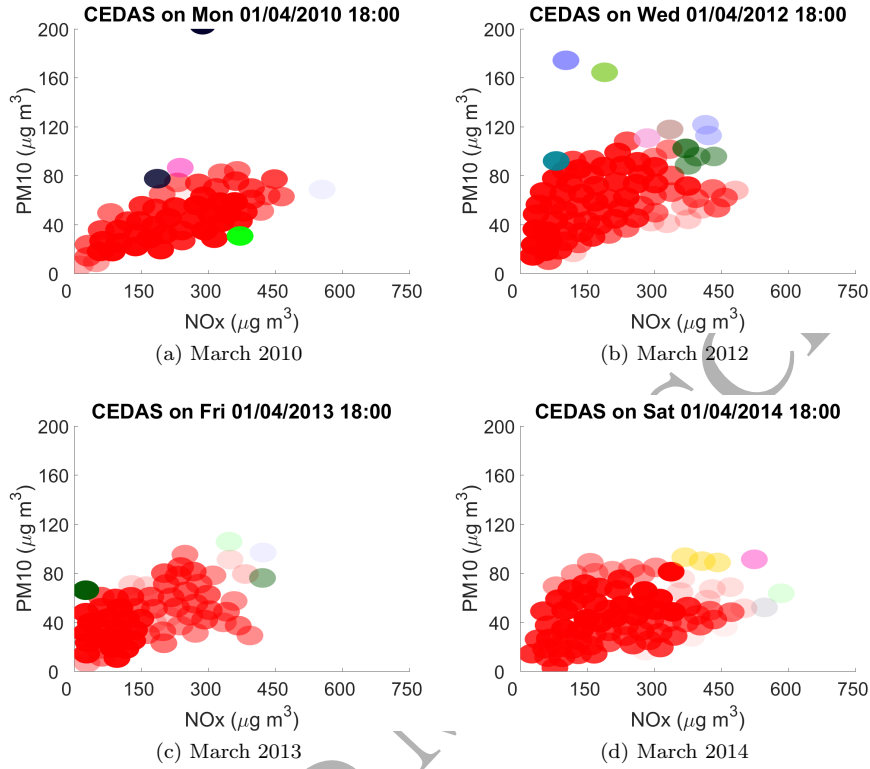


Figure 14: Plots of CEDAS clustering with a 28 day decay period showing variation of the data for March over a 5 year period.

2012, however, shows some more extreme anomalies, shown in blue and green, with particularly high  $PM_{10}$  values. Thus we see that these anomalies detected in March 2012 were not measured in any other year. This demonstrates how CEDAS may be used to analyse yearly changes, i.e. long term shifts in data and find anomalies independent of drift.

## 5. Conclusions

A new, fully online clustering technique for clustering data into arbitrarily shaped clusters is proposed. In Section 3 the algorithm has been described. We have demonstrated the ability of CEDAS to cope with evolving data streams in a fully online manner. The proposed algorithm compares favourably with similar techniques in terms of speed, cluster purity, accuracy and memory efficiency. This is especially true when the data stream is evolving constantly or there are a high number of micro-clusters. CEDAS has a linear complexity and time penalty relative to the number of data dimensions. There is also a linear time

penalty relative to the number of micro-clusters produced which, in the worst case, results in a linear relationship between the decay time and the processing time. In practice data streams that drift at such a high rate are likely to be rare and may require a different type of analysis in any case. In the case of fairly static data the number of micro-cluster will vary little and no time penalty results from an increase in decay time.

CEDAS has been demonstrated mining data streams where it proved capable of accurately detecting anomalies within the defined time periods demonstrating possible applications in network security and atmospheric science research. These results demonstrate how CEDAS could be used to automate detection across multiple dimensions that cannot be easily visualised or to present a visualization for primary interpretation by the user.

Clustering of Evolving Data streams into Arbitrary Shapes has been demonstrated to be a robust and accurate technique with linear complexity across both data stream size and data stream dimensionality. It is a fully online technique providing constant and immediate access to the clustering results.

## 6. Future Work

The work presented in this paper demonstrates the ability of the proposed technique to accurately cluster data from evolving data streams. The macro-clusters are separated by virtue of the *radius* value of the micro-clusters. Future work should consider the micro-clusters having fuzzy membership of macro-clusters as this may reduce or remove the need to separate the macro-clusters by a pre-determined value.

No numerical analysis of the macro-clusters has been considered. Proposed future work could consider quantitative methods for measuring the clustering results. Well established shape factors such as circularity, solidity or waviness etc. may provide insight into the changing relationship between clusters over time. Macro-clusters are agglomerations of micro-clusters which suggests fractal analysis [6]. Providing some measure of the location, spread, size and shape of the macro-clusters can provide information towards a quantitative assessment of the similarity and connection between the internal cluster space and difference measure to other macro-clusters.

## Acknowledgments

This work was supported by a NERC studentship through the Coordinated Airborne Studies in the Tropics (CAST) project, grants number NE/J006262/1 and NE/J006181/1.

The authors would like to thank Dr. Gary Fuller for his help selecting data from the London Air Quality Dataset.

The second author would like to acknowledge the partial support through The Royal Society grant IE141329/2014 “Novel Machine Learning Paradigms to address Big Data Stream”

## References

- [1] Aggarwal, C. C., Watson, T. J., Ctr, R., Han, J., Wang, J., & Yu, P. S. (2003). A framework for clustering evolving data streams. *Proceedings of the 29th international conference on Very large data bases*, (pp. 81–92). URL: <http://www.vldb.org/conf/2003/papers/S04P02.pdf>. doi:10.1.1.13.8650.
- [2] Angelov, P., & Zhou, X. (2008). Evolving Fuzzy-Rule Based Classifiers From Data Streams. *IEEE Transactions on Fuzzy Systems*, 16, 1462–1474. doi:10.1109/TFUZZ.2008.925904.
- [3] Babcock, B., Babu, S., Datar, M., Motwani, R., & Widom, J. (2002). Models and issues in data stream systems. *Proceedings of the twentyfirst ACM SIGMODSIGACTSIGART symposium on Principles of database systems PODS 02*, pages, 1. URL: <http://portal.acm.org/citation.cfm?doid=543613.543615>. doi:10.1145/543614.543615.
- [4] Baruah, R. D., & Angelov, P. (2013). DEC: Dynamically evolving clustering and its application to structure identification of evolving fuzzy model. *Transaction on Cybernetics*, 44, 1–16. doi:10.1109/TCYB.2013.2291234.
- [5] Bifet, A., Holmes, G., Pfahringer, B., Kranen, P., Kremer, H., Jansen, T., & Seidl, T. (2010). MOA: Massive online analysis, a framework for stream classification and clustering. *HaCDAIS 2010*, 11, 3. URL: <http://eprints.pascal-network.org/archive/00007201/>.
- [6] Botet, R., Jullien, R., & Kolb, M. (1984). Hierarchical model for irreversible kinetic cluster formation. *Physics A: Mathematical and General*, 117, 75–79. doi:10.1088/0305-4470/117/2/009.
- [7] Cao, F., Ester, M., Qian, W., & Zhou, A. (2006). Density-based clustering over an evolving data stream with noise. In . . . *Conference on Data Mining* 2 (pp. 328–339). doi:10.1145/1552303.1552307.
- [8] Chaoji, V. (2009). *Efficient Algorithms for Mining Arbitrary*. Ph.D. thesis Rensselaer Polytechnic Institute. URL: <http://www.cs.rpi.edu/~zaki/PaperDir/PhdTheses/chaoji.pdf>.
- [9] Chaoji, V., Al Hasan, M., Salem, S., & Zaki, M. J. (2008). SPARCL: Efficient and effective shape-based clustering. *Proceedings - IEEE International Conference on Data Mining, ICDM*, (pp. 93–102). doi:10.1109/ICDM.2008.73.
- [10] Dutta Baruah, R., & Angelov, P. (2012). Evolving local means method for clustering of streaming data. *IEEE International Conference on Fuzzy Systems*, (pp. 10–15). doi:10.1109/FUZZ-IEEE.2012.6251366.

- [11] Environmental Research Group, K. C. L. (2015). London Air Quality Network :: Welcome to the London Air Quality Network Data Downloads. URL: <http://www.londonair.org.uk/london/asp/datadownload.asp>.
- [12] Ester, M., Kriegel, H. P., Sander, J., & Xu, X. (1996). A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. *Second International Conference on Knowledge Discovery and Data Mining*, (pp. 226–231). URL: <http://www.aaai.org/Papers/KDD/1996/KDD96-037>. doi:10.1.1.71.1980. arXiv:10.1.1.71.1980.
- [13] Glass, L., & Mackey, M. (2010). Mackey-Glass equation. *Scholarpedia*, 5, 6908. URL: [http://www.scholarpedia.org/article/Mackey-Glass{ }equation](http://www.scholarpedia.org/article/Mackey-Glass%7B%7Dequation). doi:10.4249/scholarpedia.6908.
- [14] Hahsler, M., Arya, S., & Mount, D. (2015). Density based clustering of applications with noise (DBSCAN) and related algorithms. URL: <https://cran.r-project.org/web/packages/dbscan/index.html><http://cran.r-project.org/package=dbscan>.
- [15] Hartigan, J. A., & Wong, M. A. (1979). Algorithm AS 136: A K-Means clustering algorithm. *Applied Statistics*, 28, 100. URL: <http://www.jstor.org/stable/10.2307/2346830?origin=crossref>. doi:10.2307/2346830.
- [16] Hettich, S., & Bay, S. D. (1999). *The UCI KDD Archive*. Technical Report University of California, Department of Information and Computer Science Irvine. URL: [http://kdd.ics.uci.edu/databases/kddcup99/kddcup.data{ }10{ }percent.gz](http://kdd.ics.uci.edu/databases/kddcup99/kddcup.data%7B%7D10%7B%7Dpercent.gz).
- [17] Hyde, R. (2016). CEDAS Matlab Implementation. URL: <https://rhyde67.github.io/CEDAS/>.
- [18] Hyde, R., & Angelov, P. (2015). A new online clustering approach for data in arbitrary shaped clusters. In *2015 IEEE 2nd International Conference on Cybernetics (CYBCONF)* (pp. 228–233). Gdynia: IEEE. URL: <http://ieeexplore.ieee.org/xpl/articleDetails.jsp?arnumber=7175937>. doi:10.1109/CYBConf.2015.7175937.
- [19] Karypis, G., Han, E.-H., & Kumar, V. (1999). Chameleon: hierarchical clustering using dynamic modeling. *Computer*, 32, 68–75. doi:10.1109/2.781637.
- [20] Liu, L.-x., Guo, Y.-f., Kang, J., & Huang, H. (2009). A three-step clustering algorithm over an evolving data stream. In *2009 IEEE International Conference on Intelligent Computing and Intelligent Systems* (pp. 160–164). IEEE volume 1. URL: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5357749>. doi:10.1109/ICICISYS.2009.5357749.

- [21] Mackey, M., & Glass, L. (1977). Oscillation and chaos in physiological control systems. *Science*, 197, 287–289. URL: <http://www.sciencemag.org/content/197/4300/287.short>. doi:10.1126/science.267326.
- [22] Macqueen, J. B. (1967). Some Methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley Symposium on Math, Statistics, and Probability* (pp. 281–297). University of California Press volume 1. URL: <http://projecteuclid.org/euclid.bsmsp/1200512992>. doi:citeulike-article-id:6083430.
- [23] Norby, R. J., De Kauwe, M. G., Domingues, T. F., Duursma, R. A., Ellsworth, D. S., Goll, D. S., Lapola, D. M., Luus, K. A., MacKenzie, A. R., Medlyn, B. E., Pavlick, R., Rammig, A., Smith, B., Thomas, R., Thonicke, K., Walker, A. P., Yang, X., & Zaehle, S. (2015). Model-data synthesis for the next generation of forest free-air CO<sub>2</sub> enrichment (FACE) experiments. *The New phytologist*, 209, 17–28. URL: <http://www.ncbi.nlm.nih.gov/pubmed/26249015>. doi:10.1111/nph.13593.
- [24] Partington, K., & Cardille, J. (2013). Uncovering Dominant Land-Cover Patterns of Quebec: Representative Landscapes, Spatial Clusters, and Fences. *Land*, 2, 756–773. URL: <http://www.mdpi.com/2073-445X/2/4/756/htm>. doi:10.3390/land2040756.
- [25] PCC (2012). *2012 Air quality updating and screening assessment*. Technical Report April Plymouth City Council Plymouth. URL: [http://www.plymouth.gov.uk/air/\\_quality/\\_updating/\\_screening/\\_assessment/\\_2012.pdf](http://www.plymouth.gov.uk/air/_quality/_updating/_screening/_assessment/_2012.pdf).
- [26] Pöelitz, C., Andrienko, G., & Andrienko, N. (2010). Finding arbitrary shaped clusters with related extents in space and time. *EuroVAST 2010: International Symposium on Visual Analytics Science and Technology*, (pp. 19–25). URL: <http://diglib.eg.org/EG/DL/PE/EuroVAST/EuroVAST10/019-025.pdf.abstract.pdf;internal{%&}action=action.digitallibrary.ShowPaperAbstract>.
- [27] Ren, J., & Ma, R. (2009). Density-based data streams clustering over sliding windows. In *2009 Sixth International Conference on Fuzzy Systems and Knowledge Discovery* (pp. 248–252). IEEE volume 5. URL: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5360620>. doi:10.1109/FSKD.2009.553.
- [28] Wan, L., Ng, W. K., Dang, X. H., Yu, P. S., & Zhang, K. (2009). Density-based clustering of data streams at multiple resolutions. *ACM Transactions on Knowledge Discovery from Data*, 3, 1–28. URL: <http://dl.acm.org/citation.cfm?doid=1552303.1552307>. doi:10.1145/1552303.1552307.
- [29] Wyche, K. P., Monks, P. S., Smallbone, K. L., Hamilton, J. F., Alfarra, M. R., Rickard, a. R., McFiggans, G. B., Jenkin, M. E., Bloss, W. J., Ryan, A. C., Hewitt, C. N., & MacKenzie, A. R. (2015). Mapping gas-phase



835 organic reactivity and concomitant secondary organic aerosol formation:  
chemometric dimension reduction techniques for the deconvolution of  
complex atmospheric data sets. *Atmospheric Chemistry and Physics*, 15,  
8077–8100. URL: [http://www.atmos-chem-phys.net/15/8077/2015/](http://www.atmos-chem-phys.net/15/8077/2015/acp-15-8077-2015.html)  
840 <http://www.atmos-chem-phys.net/15/8077/2015/>. doi:10.5194/acp-15-8077-2015.

- [30] Zhou, A., Cao, F., Qian, W., & Jin, C. (2008). Tracking clusters in evolving data streams over sliding windows. *Knowledge and Information Systems*, 15, 181–214. URL: <http://link.springer.com/10.1007/s10115-007-0070-x>. doi:10.1007/s10115-007-0070-x.



845

Richard Hyde is a 3rd year PhD student researching Advanced Analysis and Visualization Techniques for Atmospheric Science and is part of the Data Science Group at Lancaster University. The research is part of the NERC sponsored Co-Ordinated Airborne Studies in the Tropics (CAST) project and aims to find  
850 new methods for improving future collection and analysis of atmospheric science data.



855

Plamen Angelov (MEng98 Sofia Technical University; PhD93 Bulgarian Academy of Sciences) holds a Chair of Intelligent Systems and leads the Data Science Group at Lancaster University, UK. He is a Senior Member of IEEE and INNS of which he is a member of the Board of Governors. He authored or  
860 co-authored over 200 peer reviewed publications including two research monographs, a dozen edited books and five patents. His research interests include autonomous machine learning, knowledge extraction from data streams, evolving systems. Prof. Angelov chairs a Technical Committee and a couple of Task Forces within IEEE and a number of high profile conferences.



Rob MacKenzie has expertise in computer simulation of atmospheric aerosol and the effects of vegetation on atmospheric composition. His work on urban sustainability more broadly includes interdisciplinary tools for assessing resilience that have been applied in Birmingham, Lancaster, London, and Milan, presented in Designing Resilient Cities, and further developed through the University of Birmingham Policy Commission on Future Urban Living. Rob has also carried responsibility for major research infrastructure throughout his career: the Geophysica high-altitude research aircraft (1996-2010) and, since November 2013, the inaugural Director of the Birmingham Institute of Forest Research. BIFoR is initiating a 10m Free-Air Carbon dioxide Enrichment (FACE) facility, one of 4 parts of a uniquely ambitious global research platform for the study of the resilience of forests under environmental change.